

Regression with Gaussian Process Networks

Manel Martínez-Ramón

ECE, UNM

October, 2018

- So far we have seen the expression of the predictive distribution in linear regression, where

$$p(\bar{f}_*|\mathbf{x}^*, \mathbf{y}, \mathbf{X}) = \mathcal{N}\left(\bar{\mathbf{w}}^\top \mathbf{x}^*, \mathbf{x}^{*\top} \mathbf{A}^{-1} \mathbf{x}^*\right)$$

and

$$\mathbf{A} = \sigma_n^{-2} \mathbf{X} \mathbf{X}^\top + \Sigma_p^{-1}$$

$$\bar{\mathbf{w}} = \left(\mathbf{X} \mathbf{X}^\top + \sigma_n^2 \Sigma_p^{-1}\right)^{-1} \mathbf{X} \mathbf{y}$$

Regression in Hilbert spaces

- Now consider a nonlinear transformation into a Hilbert Space with the form $\varphi(\mathbf{x})$.
- The corresponding expressions for the inverse of the covariance matrix \mathbf{A} and the linear parameters \mathbf{w} are:

$$\mathbf{A} = \sigma_n^{-2} \boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \boldsymbol{\Sigma}_p^{-1}$$
$$\bar{\mathbf{w}} = \left(\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \sigma_n^2 \boldsymbol{\Sigma}_p^{-1} \right)^{-1} \boldsymbol{\Phi} \mathbf{y} \quad (1)$$

where $\boldsymbol{\Phi} = \{\varphi(\mathbf{x}[1]) \cdots \varphi(\mathbf{x}[N])\}$

Note that $\boldsymbol{\Phi} \boldsymbol{\Phi}^\top$ is high dimensional, possibly infinite.

The proofs in Rasmussen et al, 2006 (page 12) are slightly different from the ones presented here.

- Now we apply the Representer Theorem in the transformation

$$\varphi(x) \rightarrow \Sigma_p^{1/2} \varphi(x)$$

- This is, the parameter vector is a linear combination of the transformed data

$$\mathbf{w} = \Sigma_p^{1/2} \Phi \alpha$$

- Also, the estimator is now

$$\bar{f}_* = \varphi(x^*)^\top \Sigma_p^{1/2} \mathbf{w}$$

- With the expressions of slide ??, the estimator can be written as

$$\bar{f}_* = \varphi(x^*)^\top \Sigma_p^{1/2} \Sigma_p^{1/2} \Phi \alpha = \varphi(x^*)^\top \Sigma_p \Phi \alpha \quad (2)$$

- Note that now the dot product between test and training data is $\varphi(x^*)^\top \Sigma_p \Phi$
- Also, note that expression $\Sigma_p \Phi \alpha$ in equation (??) plays the role of a new set of parameters, i.e.

$$\bar{f}_* = \varphi(x^*)^\top \Sigma_p \Phi \alpha = \varphi(x^*)^\top \mathbf{w}'$$

Then, we can apply to this new set of parameters the solution of equation (??) of slide ??.

- This is, for

$$\mathbf{w}' = \Sigma_p \Phi \alpha$$

we can apply the solution

$$\bar{\mathbf{w}}' = \left(\Phi \Phi^\top + \sigma_n^2 \Sigma_p^{-1} \right)^{-1} \Phi \mathbf{y}$$

- and the following equalities hold

$$\Sigma_p \Phi \alpha = \left(\Phi \Phi^\top + \sigma_n^2 \Sigma_p^{-1} \right)^{-1} \Phi \mathbf{y}$$

$$\left(\Phi \Phi^\top + \sigma_n^2 \Sigma_p^{-1} \right) \Sigma_p \Phi \alpha = \Phi \mathbf{y}$$

- In order to find α we just need to premultiply it by Φ^\top

$$\begin{aligned}\Phi^\top \left(\Phi \Phi^\top + \sigma_n^2 \Sigma_p^{-1} \right) \Sigma_p \Phi \alpha &= \Phi^\top \Phi y \\ \left(\Phi^\top \Phi \Phi^\top \Sigma_p \Phi + \sigma_n^2 \Phi^\top \Sigma_p^{-1} \Sigma_p \Phi \right) \alpha &= \Phi^\top \Phi y \\ \left(\Phi^\top \Phi \Phi^\top \Sigma_p \Phi + \sigma_n^2 \Phi^\top \Phi \right) \alpha &= \Phi^\top \Phi y \\ \left(\Phi^\top \Sigma_p \Phi + \sigma_n^2 \mathbf{I} \right) \alpha &= y\end{aligned}$$

- Finally

$$\alpha = \left(\Phi^\top \Sigma_p \Phi + \sigma_n^2 \mathbf{I} \right)^{-1} y = (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} y$$

- In expression $\boldsymbol{\alpha} = (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$ we assume that

$$\mathbf{K} = \boldsymbol{\Phi}^\top \boldsymbol{\Sigma}_p \boldsymbol{\Phi}$$

which is a matrix for which entries are

$$k(\mathbf{x}, \mathbf{z}) = \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\Sigma}_p \boldsymbol{\varphi}(\mathbf{z}) = \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\Sigma}_p^{1/2} \boldsymbol{\Sigma}_p^{1/2} \boldsymbol{\varphi}(\mathbf{z})$$

which turns out to be a kernel dot product of vectors in a Hilbert space which have been linearly transformed by matrix $\boldsymbol{\Sigma}_p^{1/2}$. As we proved in the previous chapter, this is a valid dot product.

- The parameter covariance matrix is then hidden in the dot product. This means that the choice of this matrix (which defines the prior on the parameters) is implicit in the choice of the kernel.

Outcomes of this lesson

After this lesson students should be able to

- Derive the solution for the dual parameters α of the Kernel Gaussian process for regression.
- Prove that the kernel dot product between data is

$$k(\mathbf{x}_i, \mathbf{x}_j) = \varphi^\top(\mathbf{x}_i) \Sigma_p \varphi(\mathbf{x}_j)$$